



TITLE:

# Effects of kernel function on nu support vector machines in extreme cases

AUTHOR(S):

Ikeda, K

---

CITATION:

Ikeda, K. Effects of kernel function on nu support vector machines in extreme cases. IEEE TRANSACTIONS ON NEURAL NETWORKS 2006, 17(1): 1-9

ISSUE DATE:

2006-01

URL:

<http://hdl.handle.net/2433/50322>

RIGHT:

(c)2006 IEEE. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE.

# Effects of Kernel Function on Nu Support Vector Machines in Extreme Cases

Kazushi Ikeda, *Member, IEEE*

**Abstract**—How we should choose a kernel function in support vector machines (SVMs), is an important but difficult problem. In this paper, we discuss the properties of the solution of the  $\nu$ -SVM's, a variation of SVM's, for normalized feature vectors in two extreme cases: All feature vectors are almost orthogonal and all feature vectors are almost the same. In the former case, the solution of the  $\nu$ -SVM is nearly the center of gravity of the examples given while the solution is approximated to that of the  $\nu$ -SVM with the linear kernel in the latter case. Although extreme kernels are not employed in practice, analyses are helpful to understand the effects of a kernel function on the generalization performance.

**Index Terms**—Asymptotic properties, generalization ability, kernel method,  $\nu$ -SVM, support vector machine (SVM).

## I. INTRODUCTION

IN A DECADE, support vector machines (SVMs) have attracted much attention as a new classification technique with good generalization ability [1]–[5]. The basic idea of SVMs is to map input vectors into a high-dimensional feature space and linearly separate the feature vectors with an optimal hyperplane in terms of margins, i.e., distances of given examples from a separating hyperplane. The generalization ability of SVMs has been analyzed, mainly in the framework of the PAC learning [6] where the VC dimension plays an important role [7]. Recently, studies on a more practical criterion, the average generalization error, have also been presented [8]–[12].

Another important topic regarding SVMs is how we should choose a kernel function, which has a well-defined feature space. Since any positive semidefinite function can be a kernel function, we can make a new kernel function  $K$  such as

$$K := \alpha_1 K_1 + \alpha_2 K_2 \quad (\alpha_1, \alpha_2 > 0) \quad (1)$$

$$K := K_1 K_2 \quad (2)$$

$$K := \exp[K_1] \quad (3)$$

from two arbitrary kernel functions  $K_1$  and  $K_2$  where  $:=$  means definition. Hence we need to clarify which kernel is suitable for given data. Such a problem is called “learning kernels” and has been intensively studied. For example, [13] showed how to determine hyperparameters in a set of parametric kernel functions from the viewpoint of model selection in a Bayesian framework and [14] optimized the kernel function as a problem of transduction. However, both assumed a fixed set of parametric kernel functions and reduced the problem to parameter estima-

tion. Therefore, the problem of learning kernels appears to remain open.

The purpose of this study is to contribute to this important problem by elucidating the effects of the properties of a kernel function on the SVM solutions. From this viewpoint, some asymptotic properties of SVM's with the Gaussian kernel have been reported in [15] when the two parameters of the kernel method, i.e., the steepness  $\sigma^2$  of the kernel and the softness  $C$  of the margins, go to null or infinity. Although the results are important for learning kernels, some seem strange or unusual. For example, when the parameter  $C$ , which determines the significance of constraint violation, approaches null, all the examples are classified to the same category. This is caused by the formulation that positive and negative examples are separately treated, and they are unbalanced. Hence, as an alternative, we analyzed the so-called nu support vector machines ( $\nu$ -SVM's), which are a variation of the SVM's proposed in [16] and which do not distinguish positive and negative examples from a geometrical viewpoint when homogeneous separating hyperplanes are assumed [12], [17].

In this paper, we do not restrict the kernel function to a parametric model but assume only its property in two extreme cases: One is that the diagonal elements are unity and the off-diagonal elements are almost null, that is, all feature vectors are almost orthogonal; The other is that the diagonal elements are unity and the off-diagonal elements are almost unity, that is, all feature vectors are almost the same. As a result, it is shown that the solution of the  $\nu$ -SVM is nearly the center of gravity of the given examples in the former case while the solution is approximated to that of the  $\nu$ -SVM with the linear kernel in the latter case. Although such extreme kernels are not employed in practice, analyses are helpful to understand the effects of a kernel function on the generalization performance.

The rest of the paper is organized as follows: Section II introduces the  $\nu$ -SVM and its geometrical interpretation. We analyze asymptotic properties on the solution of  $\nu$ -SVM in cases where the inner product of two distinct input vectors always takes very small or very large values; that is, null or unity in Section III. The results are applied to the cases discussed in [15] in Section IV, and are confirmed by computer simulations in Section V. Conclusions are given in Section VI.

## II. THE NU SVMs

### A. Formulation

SVMs are a kind of kernel methods; that is, they nonlinearly map an input vector  $\mathbf{x}$  to a feature vector  $\mathbf{f}(\mathbf{x})$  and separate the feature vector linearly in a high-dimensional feature

Manuscript received March 1, 2005; revised June 30, 2005. This work was supported by Grant-in-Aid for Scientific Research (14084210, 15700130) from the Ministry of Education, Culture, Sports, Science and Technology of Japan.

The author is with the Graduate School of Informatics, Kyoto University, Kyoto 606-8501, Japan (e-mail: kazushi@i.kyoto-u.ac.jp).

Digital Object Identifier 10.1109/TNN.2005.860832

space. Among many hyperplanes in the feature space that correctly separate all given examples, an SVM chooses one that maximizes the margin defined as the minimum distance of examples from a separating hyperplane. We consider a homogeneous linear dichotomy in the feature space called a Perceptron, whose separating function is represented by  $\mathbf{w}'\mathbf{f}(\mathbf{x})$  where  $\mathbf{w} \in R^M$  is called the parameter vector and  $'$  denotes the transpose. Note that an inhomogeneous linear dichotomy whose separating function is represented by  $\mathbf{w}'\mathbf{f}(\mathbf{x}) + b$  is easily transformed to a homogeneous one  $\tilde{\mathbf{w}}'\tilde{\mathbf{x}}$  by lifting-up; that is, by using augmented vectors  $\tilde{\mathbf{w}} := (\mathbf{w}; b)$  and  $\tilde{\mathbf{f}}(\mathbf{x}) := (\mathbf{f}(\mathbf{x}); 1)$  where  $(\cdot; \cdot) := (\cdot', \cdot)'$ .

According to the concept mentioned above, given  $N$  examples  $(\mathbf{x}^{(n)}, y^{(n)})$ ,  $n = 1, \dots, N$ , the  $\nu$ -SVM proposed in [16] solves the following optimization problem:

$$\min_{\mathbf{w}, \xi_n, \beta} \left[ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \beta \right] \quad \text{s.t. } \mathbf{w}'\mathbf{f}^{(n)} \geq \beta - \xi_n, \quad \xi_n \geq 0 \quad (4)$$

where  $\mathbf{f}^{(n)} := y^{(n)}\mathbf{f}(\mathbf{x}^{(n)})$  and  $C$  is a constant for soft margins [18]. Note that the original  $\nu$ -SVM introduced in [16] is formulated using inhomogeneous separating hyperplanes as is in Appendix II and is of a slightly different form to (4). However, the equivalence between the original and (4) can easily be proven taking into account  $1/\nu = CN$  as discussed in [12]. If the variable  $\beta$  is fixed to unity, (4) results in the original SVMs [1], [2].

The problem (4) is known to be equivalent to the following optimization problem called the dual problem:

$$\min_{\alpha} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t. } \mathbf{w} = \sum_{n=1}^N \alpha_n \mathbf{f}^{(n)}, \quad 0 \leq \alpha_n \leq C, \quad \sum_{n=1}^N \alpha_n = 1 \quad (5)$$

where  $\alpha_n$  are the Lagrange multipliers [12], [17].

One property of the kernel methods including the  $\nu$ -SVM's is that the so-called kernel trick is applicable; that is, the inner product of a parameter vector  $\mathbf{w}$  in (5) and a feature vector  $\mathbf{f}(\mathbf{x})$  can be calculated without the explicit expression of feature vectors as follows:

$$\mathbf{w}'\mathbf{f}(\mathbf{x}) = \sum_{n=1}^N \alpha_n \mathbf{f}^{(n)'}\mathbf{f}(\mathbf{x}) \quad (6)$$

$$= \sum_{n=1}^N \alpha_n y^{(n)} K(\mathbf{x}^{(n)}, \mathbf{x}) \quad (7)$$

where  $K(\cdot, \cdot)$  is a kernel function that determines the inner product in the feature space. In fact, Mercer's theorem shows that a nonlinear function  $\mathbf{f}(\cdot)$  exists if and only if the kernel function is positive semidefinite. Therefore, choosing a kernel function means determining a nonlinear feature map and vice versa.

Another property is that the feature vector  $\mathbf{f}(\mathbf{x}^{(n)})$  and the corresponding output  $y^{(n)}$  of an example  $\mathbf{x}^{(n)}$  always appear together in the form of  $\mathbf{f}^{(n)} = y^{(n)}\mathbf{f}(\mathbf{x}^{(n)})$ . This means

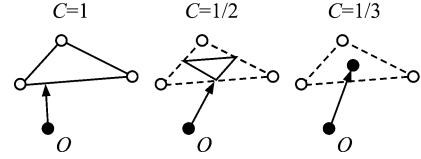


Fig. 1. Reduced convex hulls and the nearest points for  $C = 1$ ,  $C = 1/2$ , and  $C = 1/N$  when  $N = 3$ .

that an example  $(\mathbf{f}(\mathbf{x}^{(n)}), y^{(n)})$  is perfectly equivalent to  $(-\mathbf{f}(\mathbf{x}^{(n)}), -y^{(n)})$  in (5) and hence we do not have to assume the imbalance of positive and negative examples. For this reason, we call  $\mathbf{f}^{(n)}$  an example in the following. Note that if we regard all examples as chosen from one class, the formulation above is almost equivalent to the so-called one-class support vector machines [19] where the separating function is  $\mathbf{w}'\mathbf{f}(\mathbf{x}) - \beta$  instead of  $\mathbf{w}'\mathbf{f}(\mathbf{x})$ .

### B. Reduced Convex Hull

One important advantage of the  $\nu$ -SVMs in the SVM family is that (5) has a clear geometrical meaning. Minimizing the cost function  $(1/2)\|\mathbf{w}\|^2$  is equivalent to finding the point  $\hat{\mathbf{w}}$  nearest the origin that satisfies the constraints in (5). When  $C = 1$ , this restriction means that  $\mathbf{w}$  belongs to the convex hull of  $\{\mathbf{f}^{(n)}\}$ , since the sum of nonnegative weights,  $\alpha_n$ , is unity. For an arbitrary  $C < 1$ , the set in which  $\hat{\mathbf{w}}$  can exist is reduced to the so-called reduced convex hull [20], [21] by the restriction  $\alpha_n \leq C$ ; since one example can not contribute so much,  $\hat{\mathbf{w}}$  should consist of more examples, e.g., two examples when  $C = 1/2$ . The reduced convex hull shrinks into the center of gravity of all the examples when  $C = 1/N$  and vanishes when  $C < 1/N$ . So, the  $\nu$ -SVM results in the problem of finding the point nearest the origin in the reduced convex hull (Fig. 1) [12], [17]. An example  $\mathbf{f}^{(n)}$  that has a positive weight  $\alpha_n > 0$  is called a support vector. The number of support vectors is related to the generalization ability in the framework of the PAC learning as Theorem 5.2 in [1]. We denote the set of indices of support vectors by  $V$  and its complement by  $\bar{V}$ .

### C. Circumscribed Hypersphere

Suppose  $\|\mathbf{f}^{(n)}\| = 1$ , that is, feature vectors are located on a hypersphere  $S$ , and consider hard-margins' case. Then, for any example  $\mathbf{f}^{(n)}$ ,  $\hat{\mathbf{w}}$  satisfies

$$\hat{\mathbf{w}}'\mathbf{f}^{(n)} \geq \|\hat{\mathbf{w}}\|^2 \quad (8)$$

since  $\hat{\mathbf{w}}$  is the nearest point and hence

$$\|\mathbf{f}^{(n)} - \hat{\mathbf{w}}\|^2 = \|\mathbf{f}^{(n)}\|^2 + \|\hat{\mathbf{w}}\|^2 - 2\hat{\mathbf{w}}'\mathbf{f}^{(n)} \quad (9)$$

$$\leq \|\mathbf{f}^{(n)}\|^2 - \|\hat{\mathbf{w}}\|^2 = \text{const.} \quad (10)$$

Since the equality holds when  $\mathbf{f}^{(n)}$  is a support vector, the support vectors are equidistant from  $\hat{\mathbf{w}}$ ; i.e., they lie on a circumscribed hypersphere centered at  $\hat{\mathbf{w}}$ , and the other examples are inside the hypersphere [22]. Since margin maximization corresponds to radial minimization, this hypersphere is the smallest that covers all examples (see Fig. 2), [17].

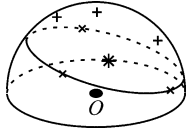


Fig. 2. The SVM solution (asterisk) corresponds to the center of the smallest ball including all examples in the feature space. Support vectors (crosses) are on the surface while the others (pluses) are inside on the unit hypersemisphere of the examples.

### III. ASYMPTOTIC PROPERTIES OF $\nu$ -SVM SOLUTIONS

We discuss the asymptotic properties of  $\nu$ -SVMs under the assumption that feature vectors are normalized,  $\|\mathbf{f}(\mathbf{x})\| = 1$ .

In the first case, the kernel function takes a very small value, nearly null, for two distinct inputs. This means that any two input vectors are almost orthogonal in the feature space. The first subsection proves that the solution of the  $\nu$ -SVM in this case is almost the gravity of center of the feature vectors.

In the second case, the kernel function takes a very large value, nearly unity due to the assumption of normalization, for two distinct input vectors. This means that any two input vectors are almost the same in the feature space. The second subsection shows that the solution of the  $\nu$ -SVM in this case is approximated to that of the so-called linear kernel SVM with inhomogeneous separating hyperplanes.

#### A. Feature Vectors Are Almost Orthogonal

Suppose that given examples are not linearly separable. The soft margin technique is one method to treat such a problem [18]; however, it sometimes reduces the generalization ability [12]. Another is to employ a kernel function that makes examples linearly separable in the feature space. For example, the Gaussian kernel has an infinite-dimensional feature space and all the feature vectors of given examples are linearly independent. This means that any set of examples becomes linearly separable in the feature space.

The condition that feature vectors are almost orthogonal is the ultimate in the latter case. Since feature vectors are little correlated to each other, it is expected that the  $\nu$ -SVM has a low generalization ability. In fact, the following theorem can be proven as below in this case:

*Theorem 1:* If the kernel function takes a small value for any pair of examples, that is

$$\left| K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \right| < \delta := \frac{1}{4N} \quad (11)$$

for any  $n \neq m$ , then any feature vector is a support vector of the  $\nu$ -SVM with hard margins.

*Proof:* Given  $N$  examples  $\mathbf{f}^{(n)}, n = 1, \dots, N$ , the solution of the  $\nu$ -SVM is the center  $\mathbf{c}_N^*$  of the minimum ball that includes all examples and is written as

$$\hat{\mathbf{w}} = \sum_{n \in V} \hat{\alpha}_n \mathbf{f}^{(n)} \quad (12)$$

where  $\hat{\alpha}_n > 0$  and  $\sum \hat{\alpha}_n = 1$ . From the discussion in Subsection II.C,  $\|\hat{\mathbf{w}} - \mathbf{f}^{(n)}\|$  for  $n \in V$  is constant and satisfies

$$\|\hat{\mathbf{w}} - \mathbf{f}^{(n)}\| > \max_{m \in \bar{V}} \|\hat{\mathbf{w}} - \mathbf{f}^{(m)}\|. \quad (13)$$

Therefore, if we show

$$\|\hat{\mathbf{w}} - \mathbf{f}^{(n)}\| < \|\hat{\mathbf{w}} - \mathbf{f}^{(m)}\| \quad (14)$$

for  $m \in \bar{V}$ , this means  $\bar{V} = \emptyset$  and the proof is completed.

Let  $\mathbf{c}$  and  $n^*$  be the center of gravity of all support vectors and the index of the farthest support vector from  $\mathbf{c}$ , respectively; that is

$$\mathbf{c} := \frac{1}{M} \sum_{n \in V} \mathbf{f}^{(n)}, \quad (15)$$

$$n^* := \arg \max_{n \in V} \|\mathbf{c} - \mathbf{f}^{(n)}\|^2 \quad (16)$$

where  $M$  is the number of support vectors,  $M := |V|$ . Then, since the ball with center  $\mathbf{c}$  and radius  $\|\mathbf{c} - \mathbf{f}^{(n^*)}\|$  includes all the support vectors, the left-hand side of (14) satisfies

$$\begin{aligned} \|\hat{\mathbf{w}} - \mathbf{f}^{(n)}\|^2 &\leq \|\mathbf{c} - \mathbf{f}^{(n^*)}\|^2 \end{aligned} \quad (17)$$

$$= \left\| \left(1 - \frac{1}{M}\right) \mathbf{f}^{(n^*)} - \sum_{n \in V - \{n^*\}} \frac{1}{M} \mathbf{f}^{(n)} \right\|^2 \quad (18)$$

$$\begin{aligned} &\leq \left(1 - \frac{1}{M}\right)^2 + \frac{M-1}{M^2} + \delta \frac{M(M-1)}{M^2} \\ &\quad + 2\delta \frac{M-1}{M} \left(1 - \frac{1}{M}\right) \end{aligned} \quad (19)$$

$$= 1 - \frac{1}{M} + \frac{\delta(M-1)(3M-2)}{M^2} \quad (20)$$

$$\leq 1 - \frac{1}{4M} \quad (21)$$

while the right-hand side satisfies

$$\begin{aligned} \|\hat{\mathbf{w}} - \mathbf{f}^{(m)}\|^2 &= 1 + \sum_{n \in V} \hat{\alpha}_n^2 - 2\mathbf{f}^{(m)} \cdot \sum_{n \in V} \hat{\alpha}_n \mathbf{f}^{(n)} + \sum_{n, l \in V, n \neq l} \hat{\alpha}_n \hat{\alpha}_l \mathbf{f}^{(n)} \cdot \mathbf{f}^{(l)} \end{aligned} \quad (22)$$

$$\begin{aligned} &\geq 1 + \sum_{n \in V} \hat{\alpha}_n^2 - 2\delta \\ &\quad - M(M-1) \left( \frac{1}{M} + 2\delta \right)^2 \delta \end{aligned} \quad (23)$$

$$\geq 1 - \frac{1}{16M}. \quad (24)$$

Here, we have used  $M \leq N$  and

$$\sum_{n \in V} \hat{\alpha}_n^2 \geq \frac{1}{M} \quad (25)$$

$$\hat{\alpha}_n \leq \frac{1}{M} + 2\delta \quad (26)$$

which are proven in Appendix I. Hence (14) holds true for any  $M$  from (21) and (24). ■

*Corollary 2:* Theorem 1 holds for the  $\nu$ -SVM with soft margins if  $C \geq (3/2N)$ .

This is proven from the facts that  $\mathbf{w}$  approaches  $\mathbf{c}$  and that  $\hat{\alpha}$  satisfies (26).

### B. Feature Vectors Are Almost the Same

Another ultimate in the opposite direction of the previous subsection is the case where the kernel function takes a very large value, nearly unity for any distinct input vectors  $\mathbf{x}_1, \mathbf{x}_2$ ,

$$K(\mathbf{x}_1, \mathbf{x}_2) \approx 1. \quad (27)$$

In this case, the set of examples  $\mathbf{f}^{(n)} = y^{(n)} \mathbf{f}(\mathbf{x}^{(n)})$  is divided into two groups; each example in one group has inner products of nearly one with the examples in the same group and those of nearly minus one with the others (see Fig. 3). That is, the property of the  $\nu$ -SVM that it does not distinguish positive and negative examples is lost here.

Since analyzing general cases is so difficult, we assume in the following that the kernel function has the form:

$$K(\mathbf{x}_1, \mathbf{x}_2) := \Psi\left(\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{\sigma^2}\right) \quad (28)$$

where  $\Psi(x)$  is a monotonically decreasing differentiable function with  $\Psi(0) = 1, \Psi(\infty) = 0$  and  $\Psi'(0) = -D_0 < 0$ , and  $\sigma$  takes a large value. One example is the Gaussian kernel where  $\Psi(x) = \exp(-x)$  and another example is the polynomial kernel with normalized inputs. Using the Taylor expansion, such a kernel function is written as

$$K(\mathbf{x}_1, \mathbf{x}_2) = K(0) - \frac{D_0}{\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_2\|^2 + \frac{D(\mathbf{x}_1, \mathbf{x}_2)}{\sigma^4} \|\mathbf{x}_1 - \mathbf{x}_2\|^4 \quad (29)$$

where

$$D(\mathbf{x}_1, \mathbf{x}_2) := \frac{1}{2} \Psi''\left(\eta \frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{\sigma^2}\right) \quad (30)$$

and  $\eta \in [0, 1]$  depending on  $\|\mathbf{x}_1 - \mathbf{x}_2\|^2$ . We set  $K(0) = 1$  and  $R = \max \|\mathbf{x}^{(n)}\|$  without loss of generality, and assume that

$$\epsilon(\mathbf{x}_1, \mathbf{x}_2) := \frac{D(\mathbf{x}_1, \mathbf{x}_2)}{\sigma^4} \|\mathbf{x}_1 - \mathbf{x}_2\|^4 \quad (31)$$

has an upper-bound  $\epsilon$  in magnitude, i.e.,  $|\epsilon(\mathbf{x}_1, \mathbf{x}_2)| \leq \epsilon$  for any input vectors  $\mathbf{x}_1, \mathbf{x}_2$ . In this case, the  $\nu$ -SVM is approximated to the inhomogeneous  $\nu$ -SVM (the original  $\nu$ -SVM) with the linear kernel  $K(\mathbf{x}_1, \mathbf{x}_2) := \mathbf{x}_1' \mathbf{x}_2$ , that is,  $\mathbf{f}(\mathbf{x}) := \mathbf{x}$ .

From the assumption (28), the cost function  $E(\alpha) = \|\mathbf{w}\|^2$  of the  $\nu$ -SVM is expressed as

$$E(\alpha) := \sum_{n,m} \alpha_n \alpha_m y^{(n)} y^{(m)} K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \quad (32)$$

$$= \left( \sum_n \alpha_n y^{(n)} \right)^2 - \sum_{n,m} \alpha_n \alpha_m y^{(n)} y^{(m)} \frac{D_0}{\sigma^2} \|\mathbf{x}^{(n)} - \mathbf{x}^{(m)}\|^2 + \sum_{n,m} \alpha_n \alpha_m y^{(n)} y^{(m)} \frac{\epsilon_{nm}}{\sigma^4} \quad (33)$$

$$= \left( \sum_n \alpha_n y^{(n)} \right)^2 + 2 \sum_{n,m} \alpha_n \alpha_m y^{(n)} y^{(m)} \frac{D_0}{\sigma^2} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} - 2 \sum_m \alpha_m y^{(m)} \sum_n \alpha_n y^{(n)} \frac{D_0}{\sigma^2} \|\mathbf{x}^{(n)}\|^2 + \sum_{n,m} \alpha_n \alpha_m y^{(n)} y^{(m)} \frac{\epsilon_{nm}}{\sigma^4} \quad (34)$$

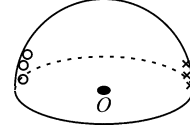


Fig. 3. Examples are divided into two groups in the feature space. One consists of positive examples (o's) and the other of negative examples (x's).

where  $\alpha := (\alpha_1, \dots, \alpha_N)$  is called a coefficient vector and  $\epsilon_{nm}$  stands for  $\epsilon(\mathbf{x}^{(n)}, \mathbf{x}^{(m)})$ . The solution of the  $\nu$ -SVM is the minimizer of  $E(\alpha)$ , denoted by  $\hat{\alpha} := (\hat{\alpha}_1, \dots, \hat{\alpha}_N)$ . Note that the last term of (33) is  $O(\sigma^{-4})$  while the second term of (33) is  $O(\sigma^{-2})$ . This means that the last term of (33) little affects the value of the cost function when  $\sigma$  takes a large value. One property of the  $\nu$ -SVM is that any example  $(\mathbf{x}^{(l)}, y^{(l)})$  s.t.  $\hat{\alpha}_l \neq 0$  satisfies

$$y^{(l)} \sum_n \hat{\alpha}_n y^{(n)} K(\mathbf{x}^{(l)}, \mathbf{x}^{(n)}) = E(\hat{\alpha}) = \sum_{n,m} \hat{\alpha}_n \hat{\alpha}_m y^{(n)} y^{(m)} K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \quad (35)$$

as mentioned in Subsection II.C.

We introduce another coefficient vector denoted by  $\alpha^P = (\alpha_1^P, \dots, \alpha_N^P)$  that satisfies  $\sum_n \alpha_n^P y^{(n)} = 0$ . Denote the set of indices for positive examples by  $V_+$  and that for negative examples by  $V_-$

$$V_+ := \{n | y^{(n)} = +1\}, \quad V_- := \{n | y^{(n)} = -1\} \quad (36)$$

and define  $\alpha^P = (\alpha_1^P, \dots, \alpha_N^P)$  as

$$\alpha_n^P := \frac{1}{1 + y^{(n)} \Delta} \hat{\alpha}_n \quad (37)$$

using the solution  $\hat{\alpha}$  of the  $\nu$ -SVM and  $\Delta := \sum_n \hat{\alpha}_n y^{(n)}$ . We can easily confirm that  $\alpha^P$  satisfies

$$\sum_{n \in V_+} \alpha_n^P = \sum_{n \in V_-} \alpha_n^P = \frac{1}{2} \quad (38)$$

and hence

$$\sum_n \alpha_n^P y^{(n)} = 0 \quad (39)$$

$$\sum_n \alpha_n^P = 1 \quad (40)$$

since

$$\sum_{n \in V_+} \hat{\alpha}_n = \frac{1 + \Delta}{2} \quad (41)$$

$$\sum_{n \in V_-} \hat{\alpha}_n = \frac{1 - \Delta}{2}. \quad (42)$$

In order to see how well  $\alpha^P$  can approximate to  $\alpha$  in separating input vectors as  $\sigma$  tends to infinity, we evaluate the difference of their outputs for an arbitrary input vector  $\mathbf{x}$ . Using the fact that  $\Delta$  is  $O(\sigma^{-2})$  in magnitude, as shown in Appendix III

$$\sum_n \hat{\alpha}_n y^{(n)} K(\mathbf{x}^{(n)}, \mathbf{x}) - \sum_n \alpha_n^P y^{(n)} K(\mathbf{x}^{(n)}, \mathbf{x}) \quad (43)$$

$$= \Delta \sum_n \frac{\hat{\alpha}_n}{1 + y^{(n)} \Delta} K(\mathbf{x}^{(n)}, \mathbf{x}) \quad (44)$$



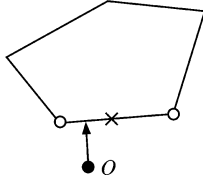


Fig. 4. Support vector (x) is expressed as a linear combination of other support vectors (o's) when they are in the same hyperplane. Here, a polygon represents the convex hull of examples.

$$\begin{aligned} &= \Delta + \Delta \frac{2D_0}{\sigma^2} \sum_n \frac{\hat{\alpha}_n}{1 + y^{(n)}\Delta} \mathbf{x}^{(n)'} \mathbf{x} \\ &\quad - \Delta \frac{D_0}{\sigma^2} \sum_n \frac{\hat{\alpha}_n}{1 + y^{(n)}\Delta} \|\mathbf{x}^{(n)}\|^2 \\ &\quad - \Delta \frac{D_0}{\sigma^2} \|\mathbf{x}\|^2 + \Delta O(\sigma^{-4}) \end{aligned} \quad (45)$$

$$= \Delta + O(\sigma^{-4}). \quad (46)$$

Hence, the difference is always  $\Delta$  with the precision of  $O(\sigma^{-2})$  and we can substitute  $\boldsymbol{\alpha}^P$  for  $\hat{\boldsymbol{\alpha}}$  by adding  $\Delta$ .

Next, we see the difference of the outputs of  $\boldsymbol{\alpha}^P$  and the solution of the inhomogeneous  $\nu$ -SVM with the linear kernel. The latter is the minimizer of

$$\begin{aligned} E^I(\boldsymbol{\alpha}) &:= \frac{2D_0}{\sigma^2} \sum_{n,m} \alpha_n \alpha_m y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} \\ &\quad \text{s.t. } 0 \leq \alpha_n \leq C, \\ &\quad \sum_{n=1}^N \alpha_n = 1, \quad \sum_n \alpha_n y^{(n)} = 0 \end{aligned} \quad (47)$$

and is denoted by  $\boldsymbol{\alpha}^I := (\alpha_1^I, \dots, \alpha_N^I)$ .

In the following, we assume the uniqueness of  $\boldsymbol{\alpha}^I$ . This assumption is made to avoid peculiar cases such that a support vector is expressed as a linear combination of other support vectors (see Fig. 4). If we remove support vectors expressed as a positively weighted sum of other examples, this is unnecessary.

Since  $|\alpha_n^P - \alpha_n^I|$  has an upper-bound of  $O(\sigma^{-1})$  as shown in Appendix IV, the difference for an arbitrary input vector  $\mathbf{x}$  is at most  $O(\sigma^{-3})$  because

$$\left| \sum_n \alpha_n^P y^{(n)} K(\mathbf{x}^{(n)}, \mathbf{x}) - \sum_n \alpha_n^I y^{(n)} K(\mathbf{x}^{(n)}, \mathbf{x}) \right| \quad (48)$$

$$\begin{aligned} &\leq \sum_n \left| (\alpha_n^P - \alpha_n^I) \frac{D_0}{\sigma^2} y^{(n)} \mathbf{x}^{(n)'} \mathbf{x} \right| \\ &\quad + \sum_n \left| (\alpha_n^P - \alpha_n^I) \frac{4\epsilon}{\sigma^4} \right| \end{aligned} \quad (49)$$

$$\leq \frac{C_2 D_0 N R^2}{\sigma^3} + \frac{4 C_2 \epsilon N}{\sigma^5} \quad (50)$$

$$= O(\sigma^{-3}). \quad (51)$$

Hence, the difference can be neglected compared to the magnitude of the output for a support vector in (35), which is  $O(\sigma^{-2})$ .

Combining both differences, it has been shown that the  $\nu$ -SVM for a large  $\sigma$  is approximated to the inhomogeneous  $\nu$ -SVM with the linear kernel by adding  $\Delta$ . Note that in this case there exists no  $\nu$ -SVM solution unless given examples are linearly separable in the input space.

#### IV. EXAMPLE: GAUSSIAN KERNEL

We apply the results in the previous section to the  $\nu$ -SVM with the Gaussian kernel and consider five asymptotic cases discussed in [15]. Note that the parameter  $C$  in the  $\nu$ -SVM is meaningful only when  $1/N \leq C \leq 1$ ; There exists no solution when  $C < 1/N$  and the solution for  $C \geq 1$  is the same as that for  $C = 1$ , whereas the conventional SVM has a unique solution for any example set and any positive parameters,  $C$  and  $\sigma^2$ . Therefore, we consider the asymptotic  $C \rightarrow 1/N$  instead of  $C \rightarrow 0$ .

*Case 1:  $\sigma^2$  Fixed and  $C \rightarrow 1/N$ :* The solution by the original SVM underfits the given examples and classifies any of them to the majority class as  $C \rightarrow 0$  since the slack variables  $\xi_n$  are ignored in the cost function [15]. The solution by the  $\nu$ -SVM, on the other hand, converges to the center of gravity of the given examples as  $C \rightarrow 1/N$ , since the reduced convex hull reduces to the point, as discussed in Subsection II-B; that is

$$\hat{\boldsymbol{\alpha}} \rightarrow \frac{1}{N} \mathbf{1} := \frac{1}{N} (1, \dots, 1). \quad (52)$$

This is a kind of underfitting since the separating hyperplane is chosen without taking into account whether the given examples are correctly classified or not.

*Case 2:  $\sigma^2$  Fixed and  $C \rightarrow \infty$ :* The solution by the original SVM classifies any of them as correctly as possible and it approaches the solution of the hard margin problem as  $C \rightarrow \infty$  [15]. The solution by the  $\nu$ -SVM is the same as that with hard margins ( $C = 1$ ) in the same way, since the reduced convex hull with  $C \geq 1$  coincides to that with  $C = 1$ .

*Case 3:  $C$  Fixed and  $\sigma^2 \rightarrow 0$ :* The solution by the original SVM overfits or underfits given examples depending on  $C$  as  $\sigma^2 \rightarrow 0$  [15]. In the  $\nu$ -SVM, however, this case satisfies the condition of Theorem 1 since  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \rightarrow \delta_{ij}$  as  $\sigma^2 \rightarrow 0$ , and hence the solution by the  $\nu$ -SVM converges to the center of gravity as discussed in Subsection III-A.

Note that we cannot simply say whether this is overfitting or underfitting because this is a kind of overfitting since the separating hyperplane is chosen so that all given examples are separated correctly, and this is a kind of underfitting since any example contributes only  $\alpha_n = 1/N$  at the same time.

*Case 4:  $C$  Fixed and  $\sigma^2 \rightarrow \infty$ :* Since  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \rightarrow 1$  as  $\sigma^2 \rightarrow \infty$ , this case satisfies the condition discussed in Subsection III-B and hence the solution of the  $\nu$ -SVM approaches that of the inhomogeneous  $\nu$ -SVM with the linear kernel.

This result is similar to that in [15]. One difference is that  $\Delta = \sum_n \hat{\alpha}_n y^{(n)}$  is always null in the original SVM while it approaches null in the  $\nu$ -SVM. Due to these properties, either of them can be approximated with an SVM with the linear kernel. Another difference is that the approximation by the linear kernel does not relate to the effect of  $C$  in the  $\nu$ -SVM while  $C$  works as if  $C/\sigma^2$  in the original SVM. This is because  $C$  appears only as the constraints in the reduced convex hull in the  $\nu$ -SVM while  $\sum_n \alpha_n$  varies according to  $C$  in the original SVM.

*Case 5:  $C \rightarrow \infty$  and  $\sigma^2 \rightarrow \infty$  With a Fixed Ratio:* Since the effect of  $C$  does not change when  $C$  becomes greater than 1, this is equivalent to the previous case in the  $\nu$ -SVM. The original SVM converges to the solution of the inhomogeneous SVM with the linear kernel with  $C/\sigma^2$  instead of  $C$ .

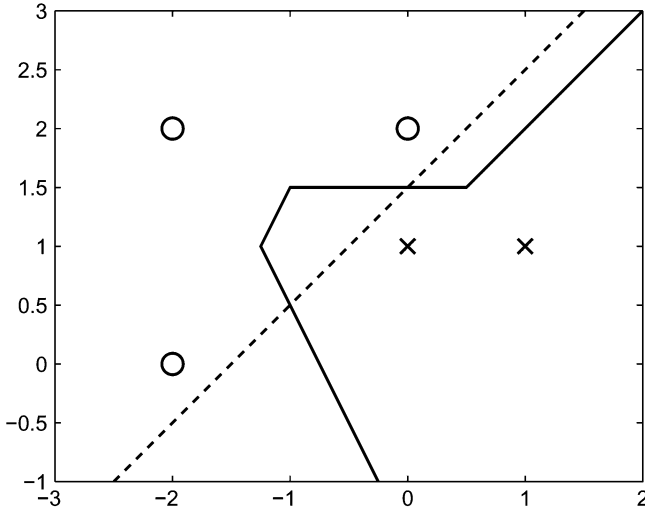


Fig. 5. Three positive examples (o's),  $\mathbf{x}^{(1)} = (0, 2)$ ,  $\mathbf{x}^{(2)} = (-2, 0)$ ,  $\mathbf{x}^{(3)} = (-2, 2)$ , and two negative examples (x's),  $\mathbf{x}^{(4)} = (0, 1)$ ,  $\mathbf{x}^{(5)} = (1, 1)$ . The solid and dashed lines express the decision boundaries for  $\sigma \rightarrow 0$  and for  $\sigma \rightarrow \infty$ , respectively.

TABLE I  
COEFFICIENT VECTORS OF  $\nu$ -SVMs WITH THE GAUSSIAN KERNEL

$\sigma^2$	$\hat{\alpha}_1$	$\hat{\alpha}_2$	$\hat{\alpha}_3$	$\hat{\alpha}_4$	$\hat{\alpha}_5$
.1	.2000	.2000	.2000	.2000	.2000
.3	.2072	.1998	.1998	.2003	.1929
1	.2814	.1779	.1728	.2375	.1303
2	.3505	.1479	.1005	.3306	.0704
3	.3854	.1454	.0438	.4019	.0236
5	.4023	.1487	0	.4489	0
10	.3925	.1383	0	.4692	0
100	.3771	.1265	0	.4963	0
1000	.3752	.1252	0	.4996	0
Center of Gravity	.2	.2	.2	.2	.2
Linear Kernel	.375	.125	0	.5	0

## V. COMPUTER SIMULATIONS

To confirm the validity of the theoretical analysis given above, some computer simulations were carried out. In each experiment, the  $\nu$ -SVM with two-dimensional (2-D) input space is given three positive examples,  $\mathbf{x}^{(1)} = (0, 2)$ ,  $\mathbf{x}^{(2)} = (-2, 0)$ ,  $\mathbf{x}^{(3)} = (-2, 2)$ , and two negative examples,  $\mathbf{x}^{(4)} = (0, 1)$ ,  $\mathbf{x}^{(5)} = (1, 1)$ , as shown in Fig. 5. The kernel function is Gaussian, that is

$$K(\mathbf{x}_1, \mathbf{x}_2) := \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{\sigma^2}\right) \quad (53)$$

and the parameter  $\sigma^2$  varies from .1 to 1000. The coefficient vector  $\hat{\alpha}$  of the  $\nu$ -SVM solution for each  $\sigma^2$  is given in Table I.

When  $\sigma^2$  is small, the vector corresponds to  $1/N$  and it approaches the coefficient vector of the inhomogeneous  $\nu$ -SVM with the linear kernel as  $\sigma^2$  increases, as proven in Section III.

## VI. CONCLUSION

We discussed the properties of the solution of the  $\nu$ -SVM in extreme cases. When all feature vectors are almost orthogonal, the solution of the  $\nu$ -SVM is nearly the center of gravity of the examples where the coefficient vector is  $1/N$ . Contrarily, when feature vectors are almost the same, the solution is approximated to that of the inhomogeneous  $\nu$ -SVM with the linear kernel. These results were confirmed by computer simulations.

These results imply that the  $\nu$ -SVM has a low generalization ability in both cases as below. Although the leave-one-out (LOO) error employed in [15] is a good approximate of the average generalization error, we cannot calculate it in the formulation of this paper. Instead, we consider in [1, Th. 5.2] here, although this holds only for the original SVM and not necessarily for the  $\nu$ -SVM with homogeneous hyperplanes. The average generalization error has three upper-bounds

$$\frac{M}{N}, \frac{R^2}{N\|\hat{\mathbf{w}}\|^2}, \frac{L}{N} \quad (54)$$

where  $L$  is the dimension of the feature space. When all feature vectors are almost orthogonal, that is, the kernel function satisfies  $|K(\cdot, \cdot)| \leq 1/(4N)$ , then

$$\frac{M}{N} = 1 \quad (55)$$

$$\frac{R^2}{N\|\hat{\mathbf{w}}\|^2} \geq \frac{4R^2}{7} \quad (56)$$

$$\frac{L}{N} \geq 1 \quad (57)$$

since

$$\|\hat{\mathbf{w}}\|^2 = \sum_{n,m} \hat{\alpha}_n \hat{\alpha}_m y^{(n)} y^{(m)} K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \quad (58)$$

$$\leq \sum_{n,m} \hat{\alpha}_n \hat{\alpha}_m \delta + (1 - \delta) \sum_n \alpha_n^2 \leq \frac{1}{4N} + \frac{3}{2N} \quad (59)$$

where we apply (64). Therefore, its generalization ability is very low. Contrarily, when feature vectors are almost the same, the generalization ability is the same as that of the inhomogeneous  $\nu$ -SVM with the linear kernel, which means that we fail to choose a kernel function in a sense.

This paper has analyzed under some restrictive assumptions what happens if we employ extreme kernel functions. Although the results are interesting from the theoretical viewpoint, they little contribute practitioners since we could expect such kernel functions would be useless. Hence, providing a more practical guideline for kernel selection remains for future work.

## APPENDIX I

### PROOF OF $\alpha_n \leq 1/M + 2\delta$

(25) is easily shown from the nonnegativeness of  $\alpha_n$ . (26) is proven as below.

Since  $\|\hat{\mathbf{w}} - \mathbf{f}^{(m)}\|$  is constant for any  $m \in V$ , so is  $\hat{\mathbf{w}} \cdot \mathbf{f}^{(m)}$ , which we denote by  $a$ . Then we have

$$\hat{\alpha}_m = -\mathbf{f}^{(m)'} \sum_{n \in V, n \neq m} \hat{\alpha}_n \mathbf{f}^{(n)} + a. \quad (60)$$

Since the inner product of distinct examples is less than  $\delta$  in magnitude and  $\sum_{n \in V} \hat{\alpha}_n = 1$ ,

$$a - \delta \leq \hat{\alpha}_m \leq a + \delta \quad (61)$$

$$\hat{\alpha}_m - \delta \leq a \leq \hat{\alpha}_m + \delta \quad (62)$$

hold true. Summing up (62) for all  $m \in V$ , we get

$$\frac{1}{M} - \delta \leq a \leq \frac{1}{M} + \delta \quad (63)$$

$$\frac{1}{M} - 2\delta \leq \hat{\alpha}_i \leq \frac{1}{M} + 2\delta \quad (64)$$

using (61) again. This means that when  $\delta$  approaches null,  $\hat{\alpha}_n$  converges to  $1/M$  and hence  $\hat{\mathbf{w}}$  to the center of gravity  $\mathbf{c}$ .

## APPENDIX II

### FORMULATION OF THE INHOMOGENEOUS $\nu$ -SVM

The original  $\nu$ -SVM's proposed in [16] employs inhomogeneous hyperplanes  $\mathbf{w}'\mathbf{f}(\mathbf{x}) + b = 0$ ; that is, hyperplanes do not necessarily include the origin. Hence, we call it the inhomogeneous  $\nu$ -SVM and formulate it as

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_n, \beta} & \left[ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \beta \right] \\ \text{s.t. } & y^{(n)} (\mathbf{w}'\mathbf{f}(\mathbf{x}^{(n)}) + b) \geq \beta - \xi_n, \quad \xi_n \geq 0. \end{aligned} \quad (65)$$

Its dual problem is easily derived to

$$\begin{aligned} \min_{\alpha_n} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t. } & \mathbf{w} = \sum_{n=1}^N \alpha_n y^{(n)} \mathbf{f}(\mathbf{x}^{(n)}), \quad 0 \leq \alpha_n \leq C, \\ & \sum_{n=1}^N \alpha_n = 1, \quad \sum_{n=1}^N y^{(n)} \alpha_n = 0 \end{aligned} \quad (66)$$

where  $\alpha_n$  are the Lagrange multipliers. Note that  $y^{(n)}$  appears alone in the last equation of (66) unlike (5).

If we define  $\tilde{\mathbf{w}} = (\mathbf{w}; b) \in \tilde{F}$  and  $\tilde{\mathbf{f}} = (\mathbf{f}; 1) \in \tilde{F}$  where  $\tilde{F}$  is the augmented feature space  $\{(\mathbf{f}, z)\} = F \times R$ , the separating hyperplane is expressed as a simple inner product  $\tilde{\mathbf{w}}'\tilde{\mathbf{f}} = 0$ ; that is, the hyperplane is homogeneous. This operation is called lifting-up (Fig. 6). Since (66) is the same as the dual problem of the homogeneous  $\nu$ -SVM in (5) except for the last equation in (66), the solution  $\hat{\mathbf{w}}$  of the inhomogeneous SVM is expressed as the point  $(\hat{\mathbf{w}}; 0)$  nearest the origin in the intersection of the hyperplane  $z = 0$  and the reduced convex hull of the lifted-up vectors  $y^{(n)}(\mathbf{f}(\mathbf{x}^{(n)}); 1)$ .

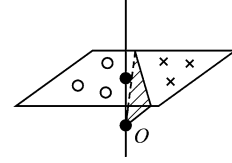


Fig. 6. Geometrical meaning of lifting-up. An inhomogeneous hyperplane (dashed line in a plane) in a space is expressed as a homogeneous hyperplane (hatched triangle) in an augmented space.

## APPENDIX III

### AN UPPER-BOUND OF $\Delta$

Summing up (35) for all  $l$  with weight  $\alpha_l^P$ , we get

$$\begin{aligned} \sum_{n,l} \hat{\alpha}_n \alpha_l^P y^{(n)} y^{(l)} K(\mathbf{x}^{(n)}, \mathbf{x}^{(l)}) \\ = \sum_{n,m} \hat{\alpha}_n \hat{\alpha}_m y^{(n)} y^{(m)} K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \end{aligned} \quad (67)$$

that is

$$\sum_{n,m} \frac{y^{(n)} \Delta}{1 + y^{(n)} \Delta} \hat{\alpha}_n \hat{\alpha}_m y^{(n)} y^{(m)} K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) = 0. \quad (68)$$

Using (29), the above equation is rewritten as

$$\begin{aligned} 0 = & \sum_{n,m} \frac{y^{(n)}}{1 + y^{(n)} \Delta} \hat{\alpha}_n \hat{\alpha}_m y^{(n)} y^{(m)} \\ & + 2 \frac{D_0}{\sigma^2} \sum_{n,m} \frac{y^{(n)}}{1 + y^{(n)} \Delta} \hat{\alpha}_n \hat{\alpha}_m y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} \\ & - \frac{D_0}{\sigma^2} \sum_{n,m} \frac{y^{(n)}}{1 + y^{(n)} \Delta} \hat{\alpha}_n \hat{\alpha}_m y^{(n)} y^{(m)} \|\mathbf{x}^{(n)}\|^2 \\ & - \frac{D_0}{\sigma^2} \sum_{n,m} \frac{y^{(n)}}{1 + y^{(n)} \Delta} \hat{\alpha}_n \hat{\alpha}_m y^{(n)} y^{(m)} \|\mathbf{x}^{(m)}\|^2 \\ & + \sum_{n,m} \frac{y^{(n)}}{1 + y^{(n)} \Delta} \hat{\alpha}_n \hat{\alpha}_m \frac{\epsilon_{nm}}{\sigma^4}. \end{aligned} \quad (69)$$

Let the first term of the right-hand side of the equation be denoted by  $T_1$ , the second by  $T_2$ , the third by  $T_3$ , the fourth by  $T_4$  and the fifth by  $T_5$ . Then, the following equations and inequalities hold true:

$$T_1 = \Delta \sum_n \frac{\hat{\alpha}_n}{1 + y^{(n)} \Delta} = \Delta \quad (70)$$

$$|T_2| \leq 2 \frac{D_0}{\sigma^2} \sum_m \hat{\alpha}_m \sum_n \frac{\hat{\alpha}_n}{1 + y^{(n)} \Delta} R^2 = 2 \frac{D_0}{\sigma^2} R^2 \quad (71)$$

$$\begin{aligned} |T_3| & \leq \frac{D_0}{\sigma^2} \left| \sum_m \hat{\alpha}_m y^{(m)} \right| \sum_n \frac{\hat{\alpha}_n}{1 + y^{(n)} \Delta} R^2 \\ & = |\Delta| \frac{D_0}{\sigma^2} R^2 \end{aligned} \quad (72)$$

$$|T_4| \leq \frac{D_0}{\sigma^2} \sum_n \frac{\hat{\alpha}_n}{1 + y^{(n)} \Delta} \sum_m \hat{\alpha}_m R^2 = \frac{D_0}{\sigma^2} R^2 \quad (73)$$



$$C_1 := \inf_{\alpha^P} \frac{\sum_{n,m} (\alpha_n^P - \alpha_n^I) (\alpha_m^P - \alpha_m^I) y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)}}{\|\alpha^P - \alpha^I\|^2}. \quad (83)$$

$$|T_5| \leq \sum_n \frac{\hat{\alpha}_n}{1 + y^{(n)} \Delta} \sum_m \hat{\alpha}_m \frac{\epsilon}{\sigma^4} = \frac{\epsilon}{\sigma^4} \quad (74)$$

which lead (68) to

$$|\Delta| \left(1 - \frac{D_0}{\sigma^2} R^2\right) \leq 3 \frac{D_0}{\sigma^2} R^2 + \frac{\epsilon}{\sigma^4}. \quad (75)$$

Therefore, the magnitude of  $\Delta$  is at most  $O(\sigma^{-2})$ .

#### APPENDIX IV AN UPPER-BOUND OF $|\alpha_n^P - \alpha_n^I|$

Using (29), an upper-bound of  $E(\alpha^P - \alpha^I)$  is given by

$$\sum_{n,m} (\alpha_n^P - \alpha_n^I) (\alpha_m^P - \alpha_m^I) y^{(n)} y^{(m)} K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \quad (76)$$

$$\leq \frac{2D_0}{\sigma^2} \sum_{n,m} (\alpha_n^P - \alpha_n^I) (\alpha_m^P - \alpha_m^I) y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} + \frac{4\epsilon}{\sigma^4} \quad (77)$$

$$\leq \frac{2D_0}{\sigma^2} \sum_{n,m} \alpha_n^P \alpha_m^P y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} - \frac{2D_0}{\sigma^2} \sum_{n,m} \alpha_n^I \alpha_m^I y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} + \frac{4\epsilon}{\sigma^4} \quad (78)$$

$$= E^I(\alpha^P) - E^I(\alpha^I) + \frac{4\epsilon}{\sigma^4} \leq \frac{6\epsilon}{\sigma^4}. \quad (79)$$

On the contrary,  $E(\alpha^P - \alpha^I)$  has a lower-bound such that

$$\sum_{n,m} (\alpha_n^P - \alpha_n^I) (\alpha_m^P - \alpha_m^I) y^{(n)} y^{(m)} K(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) \quad (80)$$

$$\geq \frac{2D_0}{\sigma^2} \sum_{n,m} (\alpha_n^P - \alpha_n^I) (\alpha_m^P - \alpha_m^I) y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} - \frac{4\epsilon}{\sigma^4} \quad (81)$$

$$= C_1 \frac{2D_0}{\sigma^2} \|\alpha^P - \alpha^I\|^2 - \frac{4\epsilon}{\sigma^4} \quad (82)$$

where (83) is shown at the top of the page. Note that  $C_1$  is positive since

$$\begin{aligned} & \sum_{n,m} (\alpha_n^P - \alpha_n^I) (\alpha_m^P - \alpha_m^I) y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} \\ & \geq \sum_{n,m} \alpha_n^P \alpha_m^P y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} \\ & - \sum_{n,m} \alpha_n^I \alpha_m^I y^{(n)} y^{(m)} \mathbf{x}^{(n)'} \mathbf{x}^{(m)} > 0 \end{aligned} \quad (84)$$

from the uniqueness of the minimizer of  $E^I(\alpha)$ . Hence, there exists a positive  $C_2$  that satisfies

$$\max_n |\alpha_n^P - \alpha_n^I| \leq \|\alpha^P - \alpha^I\| \leq \frac{C_2}{\sigma}. \quad (85)$$

#### ACKNOWLEDGMENT

The authors would like to thank T. Komoto for his assistance with computer simulations.

#### REFERENCES

- [1] V. N. Vapnik, *The Nature of Statistical Learning Theory*. New York: Springer-Verlag, 1995.
- [2] —, *Statistical Learning Theory*. New York: John Wiley and Sons, 1998.
- [3] B. Schölkopf, C. Burges, and A. J. Smola, *Advances in Kernel Methods: Support Vector Learning*. Cambridge, Cambridgeshire, UK: Cambridge Univ. Press, 1998.
- [4] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines*. Cambridge, U.K.: Cambridge Univ. Press, 2000.
- [5] A. J. Smola, P. L. Bartlett, B. Schölkopf, and D. Schuurmans, Eds., *Advances in Large Margin Classifiers*. Cambridge, MA: MIT Press, 2000.
- [6] L. G. Valiant, "A theory of the learnable," *Commun. ACM*, vol. 27, pp. 1134–1142, 1984.
- [7] V. N. Vapnik and A. Y. Chervonenkis, "On the uniform convergence of relative frequencies of events to their probabilities," *Theory Prob. Appl.*, vol. 16, pp. 264–280, 1971.
- [8] R. Dietrich, M. Oppel, and H. Sompolinsky, "Statistical mechanics of support vector networks," *Phys. Rev. Lett.*, vol. 82, no. 14, pp. 2975–2978, 1999.
- [9] M. Oppel and R. Urbanczik, "Universal learning curves of support vector machines," *Physical Review Letters*, vol. 86, no. 19, pp. 4410–4413, 2001.
- [10] K. Ikeda, "Geometry and learning curves of kernel methods with polynomial kernels," *Syst. Comput. Japan*, vol. 35, no. 7, pp. 41–48, 2004.
- [11] —, "An asymptotic statistical theory of polynomial kernel methods," *Neural Comput.*, vol. 16, no. 8, pp. 1705–1719, 2004.
- [12] K. Ikeda and T. Aoshima, "An asymptotic statistical analysis of support vector machines with soft margins," *Neural Netw.*, vol. 18, no. 3, pp. 251–259, 2005.
- [13] C. Gold and P. Sollich, "Model selection for support vector machine classification," *Neurocomputing*, vol. 55, pp. 221–249, 2003.
- [14] G. Lanckriet, N. Cristianini, P. Bartlett, L. E. Ghaoui, and M. I. Jordan, "Learning the kernel matrix with semi-definite programming," *J. Machine Learning Res.*, vol. 5, pp. 27–72, 2004.
- [15] S. S. Keerthi and C.-J. Lin, "Asymptotic behaviors of support vector machines with gaussian kernel," *Neural Comput.*, vol. 15, no. 7, pp. 1667–1689, 2003.
- [16] B. Schölkopf, A. J. Smola, R. C. Williamson, and P. L. Bartlett, "New support vector algorithms," *Neural Comput.*, vol. 12, no. 5, pp. 1207–1245, 2000.
- [17] K. Ikeda and N. Murata, "Geometrical properties of nu support vector machines with different norms," *Neural Comput.*, vol. 17, no. 11, pp. 2508–2529, 2005.
- [18] C. Cortes and V. Vapnik, "Support vector networks," *Machine Learning*, vol. 20, pp. 273–297, 1995.

- [19] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Estimating the support of a high-dimensional distribution," *Neural Comput.*, vol. 13, no. 7, pp. 1443–1471, 2001.
- [20] K. P. Bennett and E. J. Bredensteiner, "Duality and geometry in SVM classifiers," in *Proc. Int. Conf. Machine Learning*, 2000, pp. 57–64.
- [21] D. Crisp and C. Burges, "A geometric interpretation of nu-SVM classifiers," *Adv. Neural Inform. Processing Syst.*, vol. 12, pp. 244–250, 2000.
- [22] R. Herbrich, *Learning Kernel Classifiers: Theory and Algorithms*. Cambridge, MA: MIT Press, 2002.



**Kazushi Ikeda** (M'94) was born in Shizuoka, Japan, in 1966. He received the B.E., M.E., and Ph.D. degrees in mathematical engineering and information physics from The University of Tokyo, Tokyo, Japan, in 1989, 1991, and 1994, respectively.

From 1994 to 1998, he was with the Department of Electrical and Computer Engineering, Kanazawa University, Kanazawa, Japan. Since 1998, he has been with the Department of Systems Science, Graduate School of Informatics, Kyoto University, Kyoto, Japan. His research interests are focused on

adaptive and learning systems, including neural networks, adaptive filters, and machine learning.